

# Assessed coursework: Integrating R and C++

Cecina Babich Morrow

2024-01-29

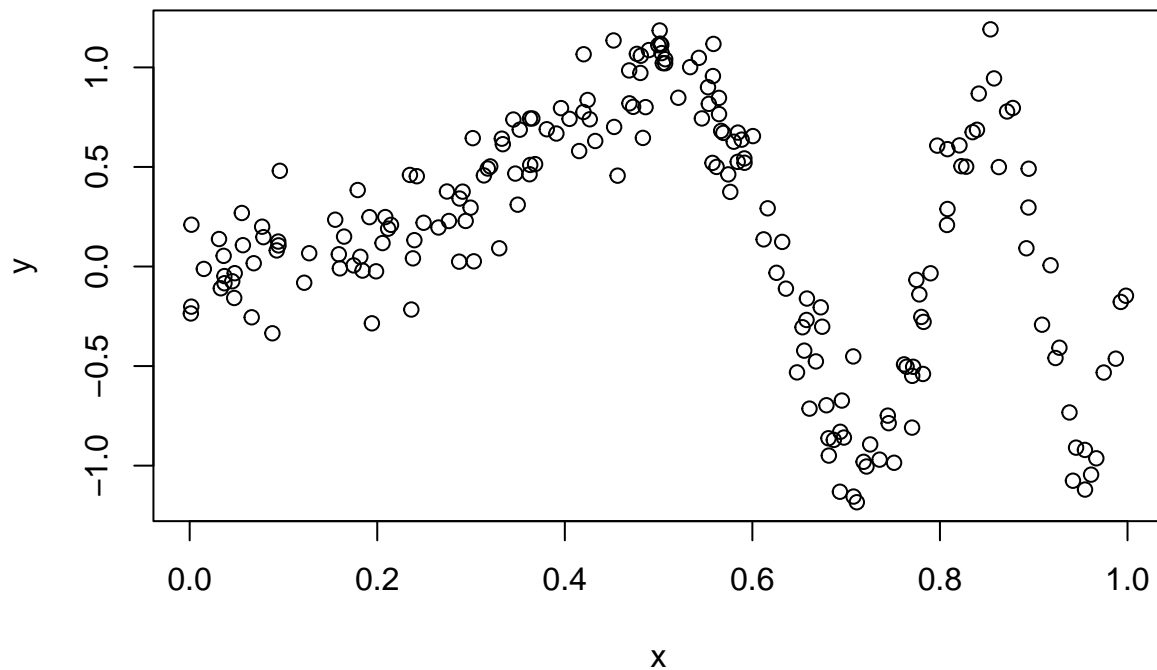
## Adaptive kernel regression smoothing

We will consider data generated from the model

$$y_i = \sin(\alpha \pi x^3) + z_i$$

where  $z_i \sim N(0, \sigma^2)$ ,  $i \in \{1, \dots, n\}$ .

```
set.seed(998)
# n = 200
nobs <- 200
x <- runif(nobs)
# alpha = 4, sigma = 0.2
y <- sin(4*pi*x^3) + rnorm(nobs, 0, 0.2)
plot(x, y)
```



We want to model this data using a kernel regression smoother (KRS) by estimating  $\mu(x) = \mathbb{E}(y|x)$ . The KRS estimator is given by

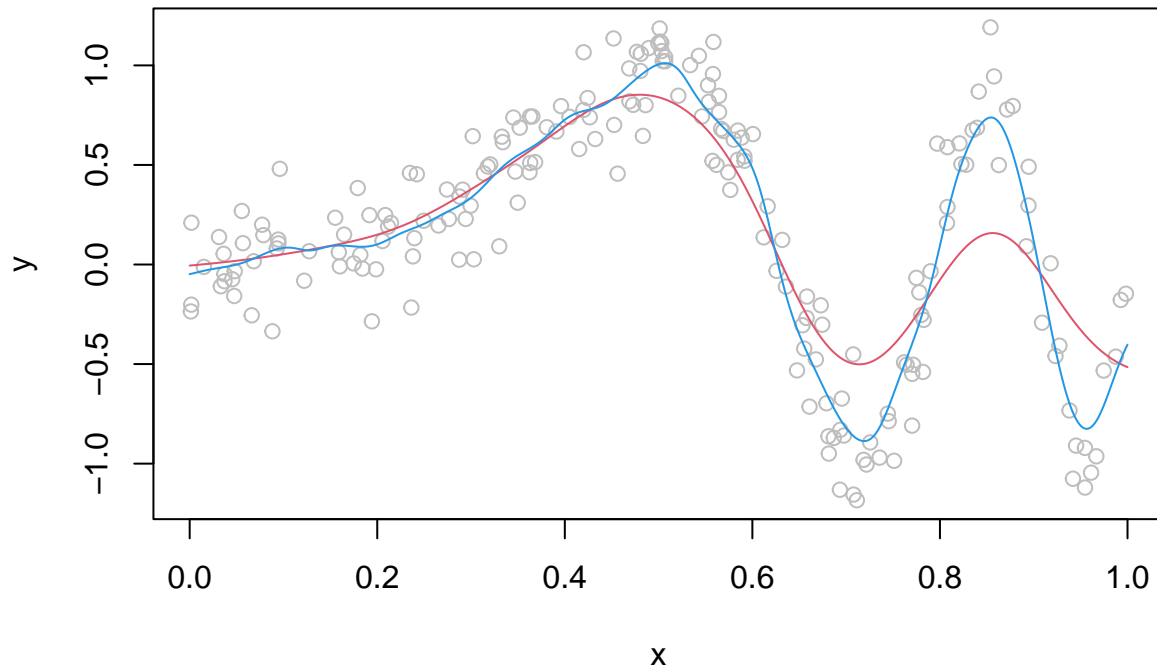
$$\hat{\mu}(x) = \frac{\sum_{i=1}^n \kappa_{\lambda}(x, x_i) y_i}{\sum_{i=1}^n \kappa_{\lambda}(x, x_i)}$$

where  $\kappa$  is a kernel function with bandwidth  $\lambda > 0$ . The following R function uses the Gaussian kernel function with variance  $\lambda^2$ :

```
meanKRS <- function(y, x, x0, lam){  
  
  n <- length(x)  
  n0 <- length(x0)  
  
  out <- numeric(n0)  
  for(ii in 1:n0){  
    out[ii] <- sum( dnorm(x, x0[ii], lam) * y ) / sum( dnorm(x, x0[ii], lam) )  
  }  
  
  return( out )  
}
```

We can compare the performance of the KRS estimator with different bandwidths:

```
xseq <- seq(0, 1, length.out = 1000)  
muSmoothLarge <- meanKRS(y = y, x = x, x0 = xseq, lam = 0.06)  
muSmoothSmall <- meanKRS(y = y, x = x, x0 = xseq, lam = 0.02)  
plot(x, y, col = "grey")  
lines(xseq, muSmoothLarge, col = 2)  
lines(xseq, muSmoothSmall, col = 4)
```



### Q1a

We want to write a C++ version of the `meanKRS` function. The function is as follows:

```
#include <R.h>
#include <Rinternals.h>
#include <Rmath.h>

SEXP meanKRS_C(SEXP y_vec, SEXP x_vec, SEXP x0_vec, SEXP lambda_param)
{
    int n = length(x_vec);
    int n0 = length(x0_vec);
    SEXP out = PROTECT(allocVector(REALSXP, n0));

    double *y = REAL(coerceVector(y_vec, REALSXP));
    double *x = REAL(coerceVector(x_vec, REALSXP));
    double *x0 = REAL(coerceVector(x0_vec, REALSXP));
    double lambda = REAL(lambda_param)[0];

    for (int i = 0; i < n0; i++)
    {
        double sum_dens_norm_y = 0;
        double sum_dens_norm = 0;
```

```

    for (int j = 0; j < n; j++)
    {
        double dens_norm = dnorm(x[j], x0[i], lambda, 0);
        sum_dens_norm_y += dens_norm * y[j];
        sum_dens_norm += dens_norm;
    }

    REAL(out)[i] = sum_dens_norm_y / sum_dens_norm;
}

UNPROTECT(1);

return out;
}

```

We can compile the code and then load the function as follows:

```

# Compile the code
system(paste0("R CMD SHLIB ", here("portfolios/02_interfacing_r_with_c++/meanKRS_C.c")))
# Load the binary code
dyn.load(here("portfolios/02_interfacing_r_with_c++/meanKRS_C.so"))
# Check if the function is loaded
is.loaded("meanKRS_C")

```

```
## [1] TRUE
```

Next, we can call it using `.Call` and compare the results with the R function:

```

c_smooth_test <- .Call("meanKRS_C",
    y_vec = y,
    x_vec = x,
    x0_vec = xseq,
    lambda_param = 0.06)

# Compare with results of R function
all.equal(muSmoothLarge, c_smooth_test)

```

```
## [1] TRUE
```

```

# Compare computing time
krs_R <- function() meanKRS(y = y, x = x, x0 = xseq, lam = 0.06)
krs_C <- function() .Call("meanKRS_C", y = y, x = x, x0 = xseq, lambda = 0.06)
library(microbenchmark)
microbenchmark(krs_R(), krs_C(), times = 500)

```

```

## Unit: milliseconds
##      expr      min       lq      mean     median        uq      max neval
##  krs_R()  9.952303 10.658409 11.226581 10.766600 11.309105 15.967211   500
##  krs_C()  2.621224  3.035951  3.150422  3.088098  3.219024  3.849808   500

```

The C version of the function is over 3.5 times faster than the R version on average.

## Q1b

We now want to implement  $k$ -fold cross-validation for selecting the bandwidth  $\lambda$ . We can first do so in R:

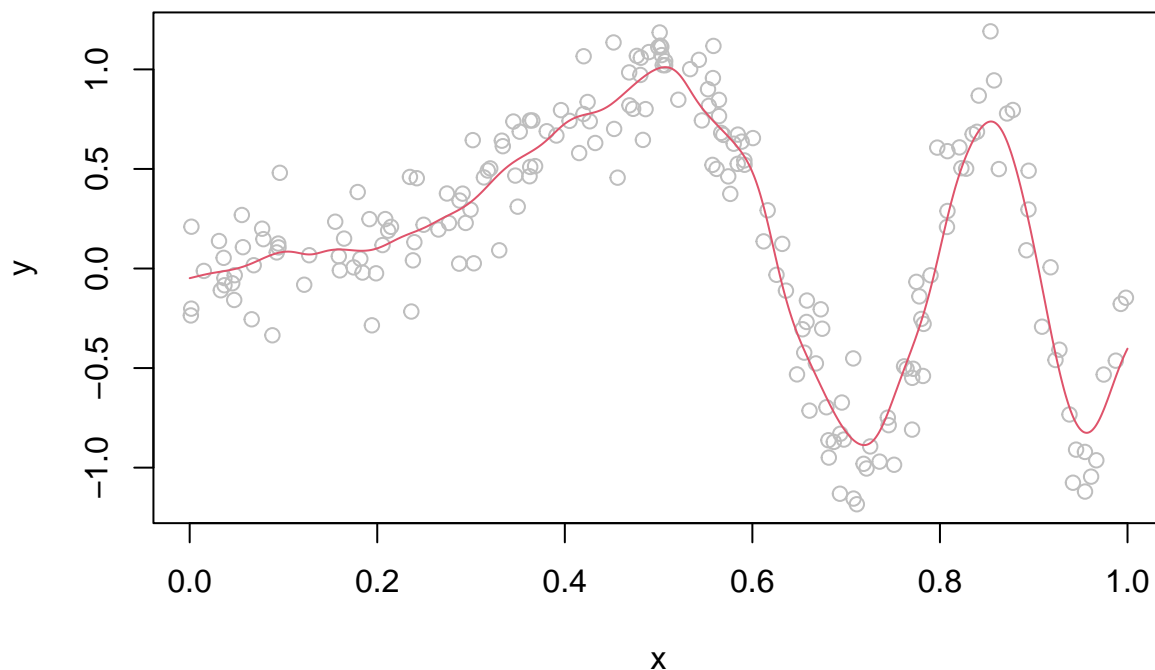
```
krsCV <- function(y, x, k, lam_seq){  
  
  n <- length(x)  
  groups <- sample(rep(1:k, length.out = n), size = n)  
  
  mse_table <- data.frame(lambda = rep(lam_seq, each = k),  
                           fold = rep(1:k, length(lam_seq)),  
                           mse = NA)  
  for (lambda in lam_seq) {  
    for (i in 1:k) {  
      # Set up training and testing sets  
      x_train <- x[groups != i]  
      y_train <- y[groups != i]  
      x_test <- x[groups == i]  
      y_test <- y[groups == i]  
  
      # Fit the model on the training set and get values for x_test  
      mu_pred <- meanKRS(y = y_train, x = x_train, x0 = x_test, lam = lambda)  
  
      # Calculate MSE on the testing set  
      mse <- mean((mu_pred - y_test)^2)  
  
      mse_table$mse[mse_table$lambda == lambda & mse_table$fold == i] <- mse  
    }  
  }  
  
  mean_mse <- mse_table %>%  
    group_by(lambda) %>%  
    summarise(mean_mse = mean(mse))  
  
  best_lambda <- mean_mse$lambda[which.min(mean_mse$mean_mse)]  
  
  return(best_lambda)  
}
```

Using this R function, we can find the value of  $\lambda$  with the lowest average mean squared error (MSE) over 5-fold cross-validation:

```
best_lambda <- krsCV(y = y, x = x, k = 5, lam_seq = seq(0.01, 0.1, by = 0.01))  
# See the best lambda value from 5-fold CV  
best_lambda
```

```
## [1] 0.02
```

```
xseq <- seq(0, 1, length.out = 1000)  
mu_best_lambda <- meanKRS(y = y, x = x, x0 = xseq, lam = best_lambda)  
plot(x, y, col = "grey")  
lines(xseq, mu_best_lambda, col = 2)
```



We can see the fit above using  $\lambda = -0.0485553, -0.0474475, -0.0463246, -0.0451891, -0.0440435, -0.0428904, -0.0417327, -0.0405804, 9.4738068 \times 10^{-4}, 0.002222, 0.0035456, 0.0049191, 0.0063434, 0.0078192, 0.0093466, 0.0109258, 0.0125564, 0.0142377, 0.0160000$

```
# dyn.unload(here("portfolios/02_interfacing_r_with_c++/krsCV_Cpp.so"))
# Compile the code
system(paste0("R CMD SHLIB ", here("portfolios/02_interfacing_r_with_c++/krsCV_Cpp.cpp")))
# Load the binary code
dyn.load(here("portfolios/02_interfacing_r_with_c++/krsCV_Cpp.so"))
# Check if the function is loaded
is.loaded("krsCV_Cpp")
```

```
## [1] TRUE
```

We can now use the C++ function to perform 5-fold cross-validation and return the value of  $\lambda$  yielding the highest average MSE across the 5 folds.

```
# Get the best lambda value from 5-fold CV
# Using the C++ function
c_best_lambda <- .Call("krsCV_Cpp",
                      y_vec = y,
                      x_vec = x,
                      k_val = as.integer(5),
                      lambda_sequence = seq(0.01, 0.1, by = 0.01))
c_best_lambda
```

```
## [1] 0.02
```

We can compare the computational performance of `krsCV_Cpp` to our R function `krsCV`:

```
# Compare computing time
krsCV_R <- function() krsCV(y = y,
                           x = x,
                           k = 5,
                           lam_seq = seq(0.01, 0.1, by = 0.01))
krsCV_Cpp <- function() .Call("krsCV_Cpp",
                              y_vec = y,
                              x_vec = x,
                              k_val = as.integer(5),
                              lambda_sequence = seq(0.01, 0.1, by = 0.01))

microbenchmark(krsCV_R(), krsCV_Cpp(), times = 500)

## Unit: milliseconds
##      expr      min      lq      mean     median        uq      max neval
## krsCV_R() 19.270213 20.537017 21.680641 21.466783 22.017895 28.495935   500
## krsCV_Cpp()  4.708349  5.233496  5.432359  5.409505  5.606198  8.414248   500
```

Once again, the C++ version is much faster than the R version, this time around 4 times faster.

## Q2

The following R function allows the bandwidth to depend on  $x$ , i.e.  $\lambda = \lambda(x)$ :

```
mean_var_KRS <- function(y, x, x0, lam){

  n <- length(x)
  n0 <- length(x0)
  mu <- res <- numeric(n)

  out <- madHat <- numeric(n0)

  for(ii in 1:n){
    mu[ii] <- sum( dnorm(x, x[ii], lam) * y ) / sum( dnorm(x, x[ii], lam) )
  }

  resAbs <- abs(y - mu)
  for(ii in 1:n0){
    madHat[ii] <- sum( dnorm(x, x0[ii], lam) * resAbs ) / sum( dnorm(x, x0[ii], lam) )
  }

  w <- 1 / madHat
  w <- w / mean(w)

  for(ii in 1:n0){
    out[ii] <- sum( dnorm(x, x0[ii], lam * w[ii]) * y ) /
               sum( dnorm(x, x0[ii], lam * w[ii]) )
  }
}
```

```

return( out )
}

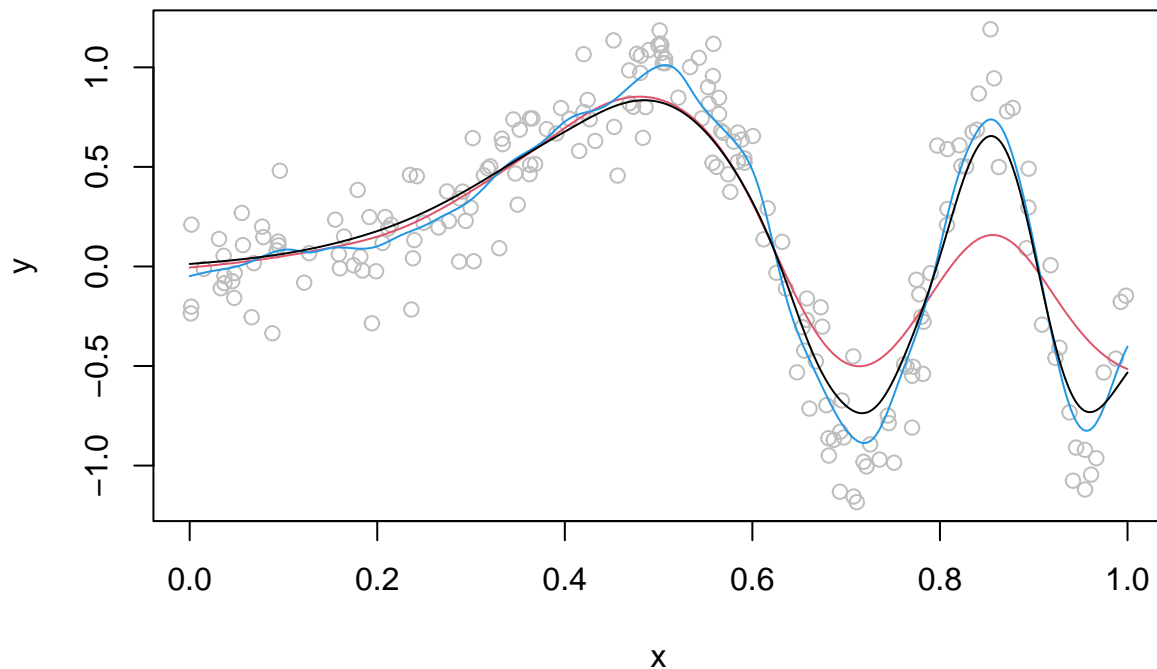
```

We can see the results here, where the red and blue lines show the fits from earlier where  $\lambda = 0.06$  and  $\lambda = 0.02$ , respectively, and the black line shows the fit using the adaptive bandwidth function `mean_var_KRS`:

```

xseq <- seq(0, 1, length.out = 1000)
muSmoothAdapt <- mean_var_KRS(y = y, x = x, x0 = xseq, lam = 0.06)
plot(x, y, col = "grey")
lines(xseq, muSmoothLarge, col = 2) # red
lines(xseq, muSmoothSmall, col = 4) # blue
lines(xseq, muSmoothAdapt, col = 1) # black

```



We now want to write a version of `mean_var_KRS` in C++ and compare the results and performance with the R function. Code for the C++ function can be found on GitHub [here](#).

```

# dyn.unload(here("portfolios/02_interfacing_r_with_c++/mean_var_krs_Cpp.so"))
# Compile the code
system(paste0("R CMD SHLIB ", here("portfolios/02_interfacing_r_with_c++/mean_var_krs_Cpp.cpp")))
# Load the binary code
dyn.load(here("portfolios/02_interfacing_r_with_c++/mean_var_krs_Cpp.so"))
# Check if the function is loaded
is.loaded("mean_var_krs_Cpp")

```

```
## [1] TRUE
```



We can use the C++ function to fit the model and compare the results with the R function:

```
mean_var_test <- .Call("mean_var_krs_Cpp",
  y_vec = y,
  x_vec = x,
  x0_vec = xseq,
  lambda_param = 0.06)
```

```
# Compare with results of R function
all.equal(muSmoothAdapt, mean_var_test)
```

```
## [1] TRUE
```

```
# Compare computing time
var_krs_R <- function() mean_var_KRS(y = y, x = x, x0 = xseq, lam = 0.06)
var_krs_C <- function() .Call("mean_var_krs_Cpp",
  y_vec = y,
  x_vec = x,
  x0_vec = xseq,
  lambda_param = 0.06)
microbenchmark(var_krs_R(), var_krs_C(), times = 500)
```

```
## Unit: milliseconds
##      expr      min       lq      mean     median        uq      max neval
## var_krs_R() 21.187948 23.813029 24.563429 24.138163 24.486816 66.225953   500
## var_krs_C()  5.909083  6.808422  6.970961  6.976432  7.055845  9.482412   500
```

Once again, the implementation in C++ is much faster than in R, a little over 3.5 times faster.